# Publish your terrestrial model data with ESS-DIVE

**Charuleka Varadharajan**

ESS-DIVE Deputy Lead, Lead: Community Engagement

# Objectives

- What are model data?
- ESS-DIVE's community-oriented process for model-data archiving
- New guidelines for storing model data with ESS-DIVE
- Other ESS-DIVE capabilities for hosting model data

- **Feedback on additional needs** for publishing model data on ESS-DIVE.

*Learn about ESS-DIVE's new guidelines for archiving model data and apply it to your datasets!*

# What does the "model data" mean to you?

# What are "model data"?

- Model outputs: Various dimensions and resolutions of final raw output files, spin-up output files, restart files, test data files, higher level output files corresponding to figures or tables
- Model input files: e.g. forcings, parameters
- Metadata files
- Model code
- Scripts for model set-up and initialization; parameterization; post-processing; and visualizations.
- Visualization files

# Model data on ESS-DIVE



Rudisill et al:
https://data.ess-dive.lbl.gov/view/doi:10.15485/1845448

# Process to develop guidelines for terrestrial model data archiving

# Review of data systems that store model data

| Data center | Storage limit per data publication | Provides data contributor guidelines | |
|---|---|---|---|
| | | Model-data specific? | Other? |
| National Science Foundation Arctic Data Center | No limit | Yes | Yes |
| Oak Ridge National Laboratory DAAC | NA[1] | Yes | Yes |
| NASA's Earth Observing System Data and Information System (EOSDIS) | NA[1] | NA[1] | Yes |
| U.S. DOE ESS-DIVE | 10GB/500 GB[2] | No | Yes |
| Dryad | 300 GB[2] | No | Yes |
| Zenodo | 50 GB | No | No |
| Earth System Grid Federation (ESGF) | NA[1] | NA[1] | NA[1] |

[1]NA: Not available, i.e. no public information found.
[2]Limit on size of individual files. For ESS-DIVE, 10GB is the default file size limit, and can be increased upto 500GB by request. Files >500GB are considered upon review.

*Simmonds et al., 2022, Data Sci. Journ.*

# Other community guidelines

NSF Earthcube Model Data Rubric/ AGU guidelines

# Community Feedback

Final Guidelines and paper

Background Research

Community Needs

Community Consensus

Case studies with published model data

Final Whitepaper

Draft White Paper

Community Feedback Webinars

ESS-DIVE

9

# Questions to Modelers

- Gathered information on
  - Models used
  - File storage and specifications
  - What is worth archiving and how long is the data useful?
  - Archiving protocols
  - Features needed from ESS-DIVE
- Circulated widely to ESS community and conducted 2 webinars on this feedback

1.1 Model data background and storage needs:

1. Name the model(s) you would like to archive data from:
2. How are your data represented spatially? (e.g., 1 km$^2$ resolution):
3. Temporal discretization and range of dataset (e.g., hourly time-step for 100 years):
4. Number of files for a typical simulation to be archived:
5. Average file size for a typical simulation to be archived:
6. Types of file formats for a typical simulation to be archived (e.g., netCDF):

1.2 What's worth archiving, why, and for how long?

1. Which model data are worth archiving? (Rate each of the following on a scale from 1-5 as Not important at all, Not so important, Somewhat important, Very important, or Extremely important)
   a. Model inputs
   b. Model outputs
   c. Metadata
   d. Model code
   e. Other scripts
   f. Model testing data
   g. Other: Describe any other type of model data worth archiving
2. In general how long do you think model data remains useful? (multiple choice)
   a. <1 year
   b. 1-2 years
   c. 5-10 years
   d. <10 years
3. Rate the importance of having these features in a model data repository on a scale of 1-5 (Not important at all, Not so important, Somewhat important, Very important, or Extremely important)
   a. Sharing of data
   b. Data preservation (for time period indicated above)
   c. Complete model data packages that can reproduce the model outputs
   d. Clear documentation
   e. Usability of archived model data
4. Describe any additional considerations that are important to informing the development of a successful model data archive.

1.3 Approaches to archiving model data

1. Does your group currently archive model data? (yes or no)
2. Describe your group's approach to archiving model data.
3. Would you be willing to share documentation of your model data archiving protocol?
4. How satisfied are you with your model data archive? (0: Very dissatisfied to 5: Very satisfied)

1.4 Last thoughts

1. Note any recommendations you have for model data storage options.
2. Would your group be willing to learn a new method for data archiving? (yes, no, absolutely not, maybe)
3. Note any other thoughts or comments here.

# Participants and projects involved

| Researcher | DOE Project(s) Represented |
| --- | --- |
| Charlie Koven | NGEE-Tropics, NGEE-Arctic |
| Jitu Kumar | NGEE-Tropics, NGEE-Arctic |
| Dipankar Dwivedi | Watershed Function SFA |
| Anthony Walker | NGEE-Tropics, FACE-MDS, Oakridge TES SFA, RUBISCO, NGEE-Arctic |
| Xingyuan Chen | PNNL SBR SFA, IDEAS, EXOSHEDS |
| Scott Painter | NGEE-Arctic, IDEAS, Oakridge Mercury SFA, Exosheds |
| Dan Ricciuto | Oakridge TES SFA |
| Qing Zhu | E3SM, RUBISCO |
| Ethan Coon | NGEE-Arctic, Exasheds |
| Maoyi Hung | NGEE Tropics, SBC SFA |
| Kate Maher | SLAC SFA |
| Ahmad Jan | NGEE-Arctic |

# Models used across DOE projects

- ELM
- ELM-FATES
- ELM-BeTR
- CLM
- PFLOTRAN
- CABLE
- SDVGM
- GDAY
- ED2
- LPJ-GUESS

- MAAT
- ATF
- ATS
- OpenFOAM (CFD)
- Crunch
- Crunch-Flow
- SWAT
- Machine Learning

# File storage and specifications

- Total files per simulation: 5 to a few million files
- Average file size: 100 MB - 2 TB
- Total current storage: 100's of MB to 100's of TB (mean = 28 TB/modeler, median = 650 GB/modeler).

| | | | | | | | | **Total annual storage needs (GB)** |
|---|---|---|---|---|---|---|---|---|
| | | | Details for typical simulation[1] to be archived | | | | | |
| **Model** | **Spatial Resolution or Representation** | **Spatial Extent** | **Temporal Resolution[2]** | **Temporal Duration** | **No. of files** | **Mean file size (GB)** | **Types of file formats** | |
| Multiple LSMs[3] | Point[4] | point | daily | 200 yrs | 300 | 0.1 | CSV | 50 |
| ELM | point | point | hourly, daily | 10 -- 20 yrs | 20 | 0.004 | netCDF | 3 |
| ELM | 1/2° -- 2° | global | monthly | 250 yrs | 2500 | 0.2 | netCDF | 15000 |
| ELM-FATES | point, ~1km, ~1degree | point, regional, and global modes | sub-daily, monthly | ~500 yrs | 1K -- 10K | 50 | netCDF | 1000 |
| FATES | point | point | <hourly | 10 yrs | 70 | 3 | netCDF | 2000 |
| ELM-PFLOTRAN | 1 -- 100m | 100m -- 10 km | hourly/daily | 10+ yrs | 10 -- 100 | 10 | HDF5, netCDF | 1000 |
| PFLOTRAN | <1m | 5-6 km | <hourly | 30 yrs | 5 | 1000 | HDF5 | 10000 |
| ATS | 100m -- 250m | 10km | daily | 10 -- 100 yrs | 20 | 100 | XML + HDF5, CSV | 1000 |
| ATS | <1 -- 100m | 10m -- 10km | daily | 10 -- 100 yrs | | 2 | XML + HDF5 | 1000 |
| ATS | 0.25m | 25m | daily | 100 yrs | 50 -- 200 | | XML + HDF5 | 10 |
| CrunchFlow | <1m | <1km | <hourly | 30 days | 100 | 0.001 | TXT | 1 |

*Simmonds et al., 2022, Data Sci. Journ.*

13

# What model data components are worth archiving?

- Journal and funding requirements for data archiving
- Archive almost everything!
- Exception outputs due to size: archive only high-level outputs corresponding to key figures/findings
- Ambiguity around model code



*Simmonds et al., 2022, Data Sci. Journ.*

\+    restart files

# How long does model data remain useful?

# Purposes for storing model data publicly

- Transparency and reproducibility
- Model intercomparisons and synthesis
- Comparison with observations
- Standards for ensuring data reusability
  - Machine-readable
  - Critical metadata
  - Workflow documents
  - 90% willing to learn new guidelines or standardized reporting format for model data

*Simmonds et al., 2022, Data Sci. Journ.*

## Importance of having these features in the repository?

[Bar chart showing importance ratings from "Not important at all" (1) to "Extremely important" (5) for: data sharing, data preservation, reproducibility, clear documentation, data usability]

# Any Questions?

# Guidelines:
# How to archive model data

**DATA SCIENCE JOURNAL**

Reading: Guidelines for Publicly Archiving Terrestrial Model Data to Enhance Usability, Intercomparis…

Share:

**Research Papers**

Guidelines for Publicly Archiving Terrestrial Model Data to Enhance Usability, Intercomparison, and Synthesis

**Authors:** Maegen B. Simmonds, William J. Riley, Deborah A. Agarwal, Xingyuan Chen, Shreyas Cholia, Robert Crystal-Ornelas, Ethan T. Coon, Dipankar Dwivedi, Valerie C. Hendrix, Maoyi Huang, Ahmad Jan, Zarine Kakalia, Jitendra Kumar, Charles D. Koven, Li Li, Mario Melara, Lavanya Ramakrishnan, Daniel M. Ricciuto, Anthony P. Walker, Wei Zhi, Qing Zhu, Charuleka Varadharajan ✉

## Abstract

Scientific communities are increasingly publishing data to evaluate, accredit, and build on published research. However, guidelines for curating data for publication are sparse for model-related research, limiting the usability of archived simulation data. In particular, there are no established guidelines for archiving data related to terrestrial models that simulate land processes and their coupled interactions with climate. Terrestrial modelers have a unique set of challenges when publishing data due to the diversity of scientific domains, research questions, and the types and scales of simulations. Researchers in the U.S. Department of Energy's (DOE) projects use a variety of multiscale models to advance robust predictions of terrestrial and subsurface ecosystem processes. Here, we synthesize archiving needs for data associated with different DOE models, and provide guidelines for publishing terrestrial model data components following FAIR (Findable, Accessible, Interoperable, Reusable) principles. The guidelines recommend archiving model inputs and testing data used in final simulation runs along with associated codes, workflow scripts, and metadata in public repositories. Researchers should consider archiving model

# Components of model data

- Metadata
- Required Data Files
  - Model Inputs
  - Model outputs
  - Model code
  - Scripts
- Optional Files
  - File-level metadata (FLMD)
  - Model Testing/Validation Data
  - Documentation or user guide

# Model Input Files

- Input files required unless publicly available elsewhere
  - Examples: climate forcings, meshes, soil parameterizations
  - Use open-sourced formats such as comma separated value (.csv) or NetCDF (.nc) formats where possible
  - File names should be unique, should only contain letters, numbers, hyphens, underscores, should not contain spaces, and should not rely on case-sensitive file systems
  - Hyperlink to specific input files in metadata and user guide

*Use external linking feature if input files are publicly available on another established repository*

# Model output files

- Includes raw and post-processed data, data supporting findings, tables, figures in a paper.
- Archive all model outputs if the size of the data files are within the repository **storage limitations (500GB/file on ESS-DIVE)**
- Use decision tree if size of model output files exceed repository storage limits
- Open-sourced formats such as comma separated value (.csv) or NetCDF (.nc) formats where possible.

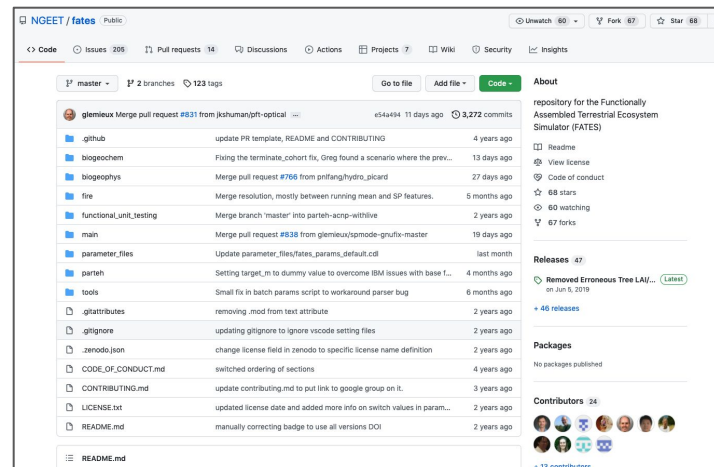*Contact* ess-dive-support@lbl.gov *if you have datasets > 0.5 TB*

# Software

- Model Code
  - Include source code(s) unless publicly available
  - Include specific version, commit hash, or citation allowing exact source code to be recovered
  - If available, include external link to tagged release in an established data repository.
  - Include links to model codes in metadata/user guide
- Scripts
  - Include **run scripts** if necessary for running model to generate results
  - Optional: Scripts for reproducing parameters and model configuration for simulations and input files, post-processing model outputs, for executing entire workflow.



*Codes published within an ESS-DIVE dataset are CC0 or CCBy4 license. Use external linking feature with software on DOE CODE*

# Optional Elements

- File-Level Metadata
- Model Testing/Validation Data
  - Data files of observations from each location in an open source format (e.g., CSV).
  - If data are publicly available in another repository, use external linking and include a reference (with DOI) in metadata and user guide.
- Documentation or user guide
  - Readme file (e.g., pdf) for each site-specific or large-scale simulation
  - Provide details on model name, version number, and required data or code dependencies.
  - Include a citation for the model code and licensing information if applicable.

# File-level metadata

**What** are file-level metadata?

- Granular information at the data file level (e.g., file name & description, start and end dates)

**Why** provide file-level metadata?

- Data users will have general understanding of info contained within a file
- FLMD can enable automatic parsing of data files so that users can eventually search & locate files across data collections

Terri Velliquette

# File-level metadata example

| File_Name | File_Description | Standard | UTC_offset |
|---|---|---|---|
| soil_samples_*.csv | 15 soil samples taken in the summer of 2019 using small hand trowel and soil probe. | csv v1.0 | - 5 hours |
| SoilPoreWaterHillslope2019.csv | 50 soil pore water samples taken from the hillslope at the site over a one year period. | EPA | - 5 hours |

- FLMD **template**: https://ess-dive.gitbook.io/file-level-metadata-reporting-format/

- Can use wildcard * to indicate when FLMD applies to multiple files

# Data package guidelines

Key Criteria

- Authorship
- File size/Repository storage capacity
- Decide what model outputs are worth archiving

*Simmonds et al., 2022, Data Sci. Journ.*

# Model data in publications

- Cite dataset and include links to the data and code publication(s) in the Data or Code Availability section and references.

# Example 1: Citing multiple datasets

*Walker, AP, et al. 2019. 'Decadal biomass increment in early secondary succession woody ecosystems is increased by CO2 enrichment'. Nature Communications, 10(1): p. 454. DOI: https://doi.org/10.1038/s41467-019-08348-1*

The site-based **meteorological dataset** (https://data.ess-dive.lbl.gov/view/ess-dive-7807cf86f1dd42a-20181127T173047368940), the **model output dataset** (https://data.ess-dive.lbl.gov/view/ess-dive-8260043c35fc925-20181130T171955541030) and the **experiment dataset** (https://data.ess-dive.lbl.gov/view/ess-dive-f525c71da7d2681-20181128T160851574946) generated and analyzed during the current study are available at the US Department of Energy's (DOE) ESS-DIVE repository.

# Example 2: Citing model code

*Koven, CD, et al. 2020. 'Benchmarking and parameter sensitivity of physiological and vegetation dynamics using the Functionally Assembled Terrestrial Ecosystem Simulator (FATES) at Barro Colorado Island, Panama'. Biogeosciences, 17(11): 3017–3044. DOI: https://doi.org/10.5194/bg-17-3017-2020*

The **FATES model** is available at https://github.com/NGEET/fates (last access: 15 May 2020; https://doi.org/10.5281/zenodo.3825474, FATES Development Team, 2020). Experiments here are based on **git commit 0bc7a5d on the fork: https://github.com/ckoven/fates (last access: 4 June 2020; https://doi.org/10.5281/zenodo.3875687**, FATES Development Team, 2019). FATES is run here within two host land surface models, CLM5 and ELMv1, available at https://github.com/ESCOMP/ctsm (git commit b9c92b7, last access: 15 May 2020; https://doi.org/10.5281/zenodo.3739617, CTSM Development Team, 2020) and https://github.com/E3SM-Project/E3SM (git commit 544db3b, last access: 15 May 2020; https://doi.org/10.11578/E3SM/dc.20180418.36, E3SM Project, 2018), respectively. Scripts to initialize parameter files and analyze model output shown here are available at https://github.com/NGEET/testbeds (last access: 15 May 2020; https://doi.org/10.5281/zenodo.3785705, Koven, 2020a), and scripts to run the all model experiments here are available at https://github.com/ckoven/runscripts (last access: 15 May 2020; https://doi.org/10.5281/zenodo.3785703, Koven, 2020b).

# Example 3: Citing data, code, scripts

*Coon, ET, et al. 2020. 'Coupling surface flow and subsurface flow in complex soil structures using mimetic finite differences'. Advances in Water Resources, 144: 103701. DOI:*
*https://doi.org/10.1016/j.advwatres.2020.103701*

The Advanced Terrestrial Simulator (ATS) (Coon et al., 2019) is open source under the BSD 3-clause license and is publicly available at https://github.com/amanzi/ats (last access: October 2019; Coon, 2016). Simulations were conducted using version 0.88. **The ATS version 0.88 is permanently stored at https://doi.org/10.5281/zenodo.3727209 (Coon et al., 2020)**. **Forcing data, model input files, Jupyter notebooks** used to generate figures, and meshes along with Jupyter notebooks used to generate the meshes are publicly available at https://doi.org/10.5440/1545603 (Jan et al., 2019). **Data products used in the model comparisons are publicly available through the NGEE Arctic long-term data archive** https://doi.org/10.5440/1416559. The observed water level can be accessed at https://doi.org/10.5440/1183767 (Liljedahl and Wilson., 2016), the soil temperature data at https://doi.org/10.5440/1126515 (Romanovsky et al., 2017), and the evapotranspiration data at https://doi.org/10.5440/1362279 (Dengel et al., 2019, respectively.
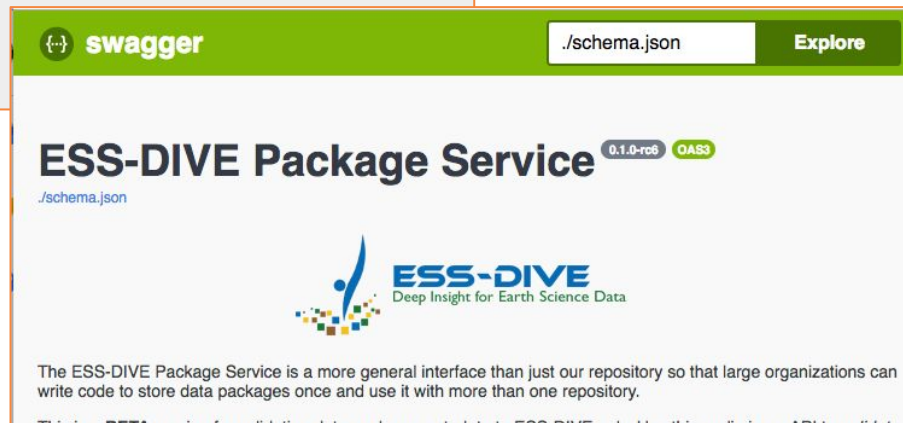
# Any Questions?

# ESS-DIVE features and plans to support model data

# Application Programming Interface (API) Upload

```
provider_spruce = {
    "name": "SPRUCE",
    "member": {
        "@id": "http://orcid.org/0000-0001-7293-3561 ",
        "givenName": "Paul J",
        "familyName": "Hanson",
        "email": "hansonpj@ornl.gov",
        "jobTitle": "Principal Investigator"
        }
    }
```



Web UI limited to 10GB/file, API allows
500GB/file

# "External Linking" Feature

Data stored on other repositories are directly linked to dataset metadata on ESS-DIVE and **some/all data do not have to be on ESS-DIVE**

Goulden T ; Hulslander D ; Hass B ; Brodie E ; Chadwick D K ; Falco N ; Maher K ; Wainwright H ; Williams K (2020): NEON AOP Imaging Spectroscopy Survey of Upper East River Colorado Watersheds: Raw-Space Radiance and Observational Variable Dataset. Watershed Function SFA, ESS-DIVE repository. Dataset. doi:el-usecase-1.2 accessed via https://data-dev.ess-dive.lbl.gov/datasets/doi:el-usecase-1.2 on 2022-03-18

# Supporting Growth In Data

- Increasing need to support large files and datasets
- Enabling multiple ways to make large data uploads seamless
  - ESS-DIVE API
  - Native ESS-DIVE Uploads
  - Globus
- Upcoming: **Tier 2 Storage layer** to for very large datasets
  - Supports very large datasets and files
  - Hierarchical data and folders
  - Direct web UI access to browse and download from Tier 2 storage
  - Globus for High Performance Downloads





Globus interface for large scale managed transfers

# Upcoming: Tier 2 Storage and Large Uploads

**Individual Files over 500GB**

Datasets containing **any file over 500GB**, such as LIDAR or drone data

**Over 100 files outside of Zip file**

Datasets containing **over 100 files** that are not stored in a **compressed (or "zipped") hierarchy** should be treated as large data.

Tier 2 will support uploading large volumes of data; large numbers of files; nested folder structures.

Globus provides high performance interface for data uploads

# Tracking Data Versions

- ESS-DIVE tracks all versions of data packages internally *before and after publication*
- Large model data may present challenges with making copies of files with different versions
- Citations indicate version used:

Creator (Publication Year). Title. Publisher. Dataset. Identifier "accessed via data.ess-dive.lbl.gov on YYYY-MM-DD"

# Any Questions?

**\*Feedback\* Tell us what you need!**

# What are your model data archiving needs?

# ESS Modeling Community Needs Summary

- Data are extremely heterogeneous and increasing in volume and complexity
- Disconnect between model and observational data
- Workflows involving manual retrieval of data are not scalable
- Short-term need for most researchers: guidelines to archive model data - e.g. those associated with publications
- Long-term: Development of a standardized *model-to-archive* pipeline and *data-to-model* pipeline.

# Future design based on community needs

- Model-to-archive pipeline
  - Community-informed guidelines on creating standardized model data packages
  - Pathway for model data packages >500 GB size threshold
  - Ability to extract specific subsets of model simulations
  - Project portals for sharing and collaborating on pre-published model data
  - (long-term) Automating the writing and/or organization of files comprising data packages for specific models or journals.
- Data-to-model pipeline
  - Support for data formats typically used in model simulations (e.g., netCDF)
  - Interoperability between individual data packages in ESS-DIVE and other data centers for model-data integration (field observations and remote sensing data to use for model development, parameterization, and performance testing to improve future measurement designs)

# Summary

**ESS-DIVE stores model data and has new terrestrial model archiving guidelines**

- Model data has many components - inputs, outputs, code, documentation

- Publishing model data is beneficial and enables reuse of simulation data

- Guidelines what to archive and how to split datasets for diverse terrestrial model data used by ESS researchers

- ESS-DIVE has many features supporting model data archiving

*Questions? Email ess-dive-support@lbl.gov*

# You can find all this material on ESS-DIVE's…



**Model Data Archiving Documentation page**
https://docs.ess-dive.lbl.gov/

**Webinar page**
where this slide deck is available for download
https://ess-dive.lbl.gov/

# Additional Resources

**Docs:** https://ess-dive.gitbook.io/model-data-archiving-guidelines/

**Paper:** https://datascience.codata.org/articles/10.5334/dsj-2022-003/

**Dataset:** https://data.ess-dive.lbl.gov/view/doi:10.15485/1813868

**Github:** https://github.com/ess-dive-community/essdive-model-data-archiving-guidelines

# Thank You!

@ESS-DIVE

**Join ESS-DIVE's Community Mailing List!**

http://bit.ly/essdiveMailingList

**Contact us at ess-dive-support@lbl.gov**